# Class 16 Mini Project

Pierce Ford (PID: A59010464)

11/19/2021

## Project Outline

1. Data Import
2. PCA
3. DESeq Analysis
4. Volcano Plot
5. Annotation
6. Pathway Analysis

## 1. Data Import

```r
metaFile <- "GSE37704_metadata.csv"
countFile <- "GSE37704_featurecounts.csv"

#Import and view metadata
colData = read.csv(metaFile, row.names=1)
head(colData)
```

```
##                 condition
## SRR493366 control_sirna
## SRR493367 control_sirna
## SRR493368 control_sirna
## SRR493369       hoxa1_kd
## SRR493370       hoxa1_kd
## SRR493371       hoxa1_kd
```

```r
#Import and view countdata
countData = read.csv(countFile, row.names=1)
head(countData)
```

```
##                 length SRR493366 SRR493367 SRR493368 SRR493369 SRR493370
## ENSG00000186092    918         0         0         0         0         0
## ENSG00000279928    718         0         0         0         0         0
## ENSG00000279457   1982        23        28        29        29        28
## ENSG00000278566    939         0         0         0         0         0
## ENSG00000273547    939         0         0         0         0         0
## ENSG00000187634   3214       124       123       205       207       212
##                 SRR493371
```

```
## ENSG00000186092        0
## ENSG00000279928        0
## ENSG00000279457       46
## ENSG00000278566        0
## ENSG00000273547        0
## ENSG00000187634      258
```

```r
# Note we need to remove the odd first $length col
countData <- as.matrix(countData[,-1])
head(countData)
```

```
##                 SRR493366 SRR493367 SRR493368 SRR493369 SRR493370 SRR493371
## ENSG00000186092         0         0         0         0         0         0
## ENSG00000279928         0         0         0         0         0         0
## ENSG00000279457        23        28        29        29        28        46
## ENSG00000278566         0         0         0         0         0         0
## ENSG00000273547         0         0         0         0         0         0
## ENSG00000187634       124       123       205       207       212       258
```

Next, let's remove rows that are all zeros.

```r
#Find empty rows
zero.rows <- which(rowSums(countData)==0)
#Remove empty rows and check that it worked
countData.filtered <- countData[-zero.rows,]
head(countData.filtered)
```

```
##                 SRR493366 SRR493367 SRR493368 SRR493369 SRR493370 SRR493371
## ENSG00000279457        23        28        29        29        28        46
## ENSG00000187634       124       123       205       207       212       258
## ENSG00000188976      1637      1831      2383      1226      1326      1504
## ENSG00000187961       120       153       180       236       255       357
## ENSG00000187583        24        48        65        44        48        64
## ENSG00000187642         4         9        16        14        16        16
```

```r
#How many genes are left?
nrow(countData.filtered)
```

```
## [1] 15975
```

## PCA

Let's check that the treated and controls cluster separately.

```r
#Generate PCA
countPCA <- prcomp(t(countData.filtered))
#Plot PCA colored by condition (knockdown or not)
plot(countPCA$x, pch=16, col=as.factor(colData$condition))
text(countPCA$x, labels=colData$condition)
```

Hooray! The clustering looks correct!

## DESeq Analysis

```
library(DESeq2)
```

```
## Loading required package: S4Vectors
```

```
## Loading required package: stats4
```

```
## Loading required package: BiocGenerics
```

```
## Loading required package: parallel
```

```
##
## Attaching package: 'BiocGenerics'
```

```
## The following objects are masked from 'package:parallel':
##
##     clusterApply, clusterApplyLB, clusterCall, clusterEvalQ,
##     clusterExport, clusterMap, parApply, parCapply, parLapply,
##     parLapplyLB, parRapply, parSapply, parSapplyLB
```

```
## The following objects are masked from 'package:stats':
##
##     IQR, mad, sd, var, xtabs


## The following objects are masked from 'package:base':
##
##     anyDuplicated, append, as.data.frame, basename, cbind, colnames,
##     dirname, do.call, duplicated, eval, evalq, Filter, Find, get, grep,
##     grepl, intersect, is.unsorted, lapply, Map, mapply, match, mget,
##     order, paste, pmax, pmax.int, pmin, pmin.int, Position, rank,
##     rbind, Reduce, rownames, sapply, setdiff, sort, table, tapply,
##     union, unique, unsplit, which.max, which.min


##
## Attaching package: 'S4Vectors'


## The following objects are masked from 'package:base':
##
##     expand.grid, I, unname


## Loading required package: IRanges


## Loading required package: GenomicRanges


## Loading required package: GenomeInfoDb


## Loading required package: SummarizedExperiment


## Loading required package: MatrixGenerics


## Loading required package: matrixStats


##
## Attaching package: 'MatrixGenerics'


## The following objects are masked from 'package:matrixStats':
##
##     colAlls, colAnyNAs, colAnys, colAvgsPerRowSet, colCollapse,
##     colCounts, colCummaxs, colCummins, colCumprods, colCumsums,
##     colDiffs, colIQRDiffs, colIQRs, colLogSumExps, colMadDiffs,
##     colMads, colMaxs, colMeans2, colMedians, colMins, colOrderStats,
##     colProds, colQuantiles, colRanges, colRanks, colSdDiffs, colSds,
##     colSums2, colTabulates, colVarDiffs, colVars, colWeightedMads,
##     colWeightedMeans, colWeightedMedians, colWeightedSds,
##     colWeightedVars, rowAlls, rowAnyNAs, rowAnys, rowAvgsPerColSet,
##     rowCollapse, rowCounts, rowCummaxs, rowCummins, rowCumprods,
##     rowCumsums, rowDiffs, rowIQRDiffs, rowIQRs, rowLogSumExps,
##     rowMadDiffs, rowMads, rowMaxs, rowMeans2, rowMedians, rowMins,
##     rowOrderStats, rowProds, rowQuantiles, rowRanges, rowRanks,
##     rowSdDiffs, rowSds, rowSums2, rowTabulates, rowVarDiffs, rowVars,
##     rowWeightedMads, rowWeightedMeans, rowWeightedMedians,
##     rowWeightedSds, rowWeightedVars
```

```
## Loading required package: Biobase

## Welcome to Bioconductor
##
##      Vignettes contain introductory material; view with
##      'browseVignettes()'. To cite Bioconductor, see
##      'citation("Biobase")', and for packages 'citation("pkgname")'.

##
## Attaching package: 'Biobase'

## The following object is masked from 'package:MatrixGenerics':
##
##      rowMedians

## The following objects are masked from 'package:matrixStats':
##
##      anyMissing, rowMedians
```

```r
#Run DESeq
dds = DESeqDataSetFromMatrix(countData=countData.filtered,
                             colData=colData,
                             design=~condition)
```

```
## Warning in DESeqDataSet(se, design = design, ignoreRank): some variables in
## design formula are characters, converting to factors
```

```r
dds = DESeq(dds)
```

```
## estimating size factors

## estimating dispersions

## gene-wise dispersion estimates

## mean-dispersion relationship

## final dispersion estimates

## fitting model and testing
```

```r
#View dds and get results
dds
```

```
## class: DESeqDataSet
## dim: 15975 6
## metadata(1): version
## assays(4): counts mu H cooks
## rownames(15975): ENSG00000279457 ENSG00000187634 ... ENSG00000276345
##    ENSG00000271254
## rowData names(22): baseMean baseVar ... deviance maxCooks
## colnames(6): SRR493366 SRR493367 ... SRR493370 SRR493371
## colData names(2): condition sizeFactor
```

```
res = results(dds)
head(res)
```

```
## log2 fold change (MLE): condition hoxa1 kd vs control sirna
## Wald test p-value: condition hoxa1 kd vs control sirna
## DataFrame with 6 rows and 6 columns
##                   baseMean log2FoldChange      lfcSE        stat      pvalue
##                  <numeric>      <numeric>  <numeric>   <numeric>   <numeric>
## ENSG00000279457    29.9136      0.1792571  0.3248216    0.551863  5.81042e-01
## ENSG00000187634   183.2296      0.4264571  0.1402658    3.040350  2.36304e-03
## ENSG00000188976  1651.1881     -0.6927205  0.0548465  -12.630158  1.43990e-36
## ENSG00000187961   209.6379      0.7297556  0.1318599    5.534326  3.12428e-08
## ENSG00000187583    47.2551      0.0405765  0.2718928    0.149237  8.81366e-01
## ENSG00000187642    11.9798      0.5428105  0.5215598    1.040744  2.97994e-01
##                        padj
##                   <numeric>
## ENSG00000279457  6.86555e-01
## ENSG00000187634  5.15718e-03
## ENSG00000188976  1.76549e-35
## ENSG00000187961  1.13413e-07
## ENSG00000187583  9.19031e-01
## ENSG00000187642  4.03379e-01
```

```
summary(res)
```

```
##
## out of 15975 with nonzero total read count
## adjusted p-value < 0.1
## LFC > 0 (up)       : 4349, 27%
## LFC < 0 (down)     : 4396, 28%
## outliers [1]       : 0, 0%
## low counts [2]     : 1237, 7.7%
## (mean count < 0)
## [1] see 'cooksCutoff' argument of ?results
## [2] see 'independentFiltering' argument of ?results
```

## Volcano plot

```
#Preliminary (i.e. boring) volcano plot
plot(res$log2FoldChange, -log(res$padj))
```

Let's improve the plot to make it more informative.

```r
#Make a  baseline color vector for all genes (will replace gray with actual color later)
mycols <- rep("gray", nrow(res))

#Color red the genes with absolute fold change above 2
mycols[abs(res$log2FoldChange) > 2] <- "red"

#Color blue those with adjusted p-value less than 0.01
#and absolute fold change more than 2
inds <- (res$padj < 0.01) & (abs(res$log2FoldChange) > 2 )
mycols[inds] <- "blue"

plot(res$log2FoldChange, -log(res$padj), col=mycols, xlab="Log2(FoldChange)", ylab="-Log(P-value)" )
```

## Annotation

```
library("AnnotationDbi")
library("org.Hs.eg.db")
```

```
##
```

```
columns(org.Hs.eg.db)
```

```
##  [1] "ACCNUM"       "ALIAS"        "ENSEMBL"      "ENSEMBLPROT"   "ENSEMBLTRANS"
##  [6] "ENTREZID"     "ENZYME"       "EVIDENCE"     "EVIDENCEALL"   "GENENAME"
## [11] "GENETYPE"     "GO"           "GOALL"        "IPI"           "MAP"
## [16] "OMIM"         "ONTOLOGY"     "ONTOLOGYALL"  "PATH"          "PFAM"
## [21] "PMID"         "PROSITE"      "REFSEQ"       "SYMBOL"        "UCSCKG"
## [26] "UNIPROT"
```

```
res$symbol = mapIds(org.Hs.eg.db,
                    keys=row.names(res),
                    keytype="ENSEMBL",
                    column="SYMBOL",
                    multiVals="first")
```

```
## 'select()' returned 1:many mapping between keys and columns
```

```r
res$entrez = mapIds(org.Hs.eg.db,
                    keys=row.names(res),
                    keytype="ENSEMBL",
                    column="ENTREZID",
                    multiVals="first")
```

```
## 'select()' returned 1:many mapping between keys and columns
```

```r
res$name =   mapIds(org.Hs.eg.db,
                    keys=row.names(res),
                    keytype="ENSEMBL",
                    column="GENENAME",
                    multiVals="first")
```

```
## 'select()' returned 1:many mapping between keys and columns
```

```r
#Check that the annotations were appended
head(res, 10)
```
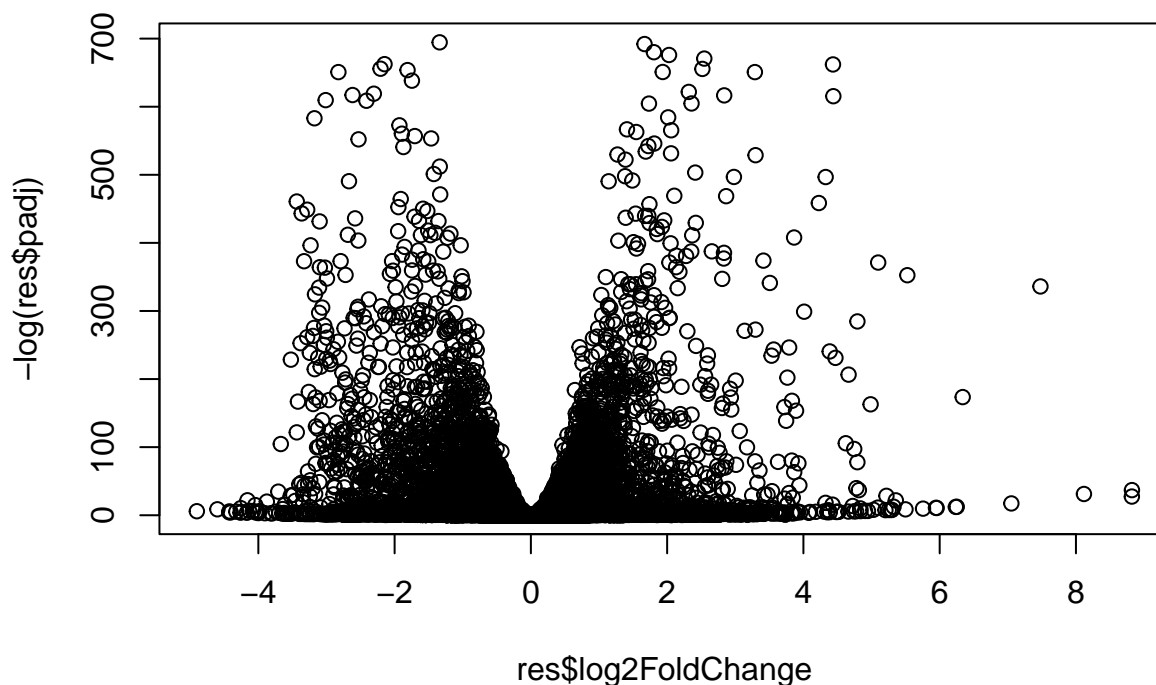
```
## log2 fold change (MLE): condition hoxa1 kd vs control sirna
## Wald test p-value: condition hoxa1 kd vs control sirna
## DataFrame with 10 rows and 9 columns
##                     baseMean log2FoldChange      lfcSE        stat      pvalue
##                    <numeric>      <numeric>  <numeric>   <numeric>   <numeric>
## ENSG00000279457    29.913579      0.1792571  0.3248216    0.551863 5.81042e-01
## ENSG00000187634   183.229650      0.4264571  0.1402658    3.040350 2.36304e-03
## ENSG00000188976  1651.188076     -0.6927205  0.0548465  -12.630158 1.43990e-36
## ENSG00000187961   209.637938      0.7297556  0.1318599    5.534326 3.12428e-08
## ENSG00000187583    47.255123      0.0405765  0.2718928    0.149237 8.81366e-01
## ENSG00000187642    11.979750      0.5428105  0.5215598    1.040744 2.97994e-01
## ENSG00000188290   108.922128      2.0570638  0.1969053   10.446970 1.51282e-25
## ENSG00000187608   350.716868      0.2573837  0.1027266    2.505522 1.22271e-02
## ENSG00000188157  9128.439422      0.3899088  0.0467163    8.346304 7.04321e-17
## ENSG00000237330     0.158192      0.7859552  4.0804729    0.192614 8.47261e-01
##                        padj      symbol      entrez                      name
##                   <numeric> <character> <character>               <character>
## ENSG00000279457 6.86555e-01       WASH9P   102723897 WAS protein family h..
## ENSG00000187634 5.15718e-03       SAMD11      148398 sterile alpha motif ..
## ENSG00000188976 1.76549e-35        NOC2L       26155 NOC2 like nucleolar ..
## ENSG00000187961 1.13413e-07       KLHL17      339451 kelch like family me..
## ENSG00000187583 9.19031e-01      PLEKHN1       84069 pleckstrin homology ..
## ENSG00000187642 4.03379e-01        PERM1       84808 PPARGC1 and ESRR ind..
## ENSG00000188290 1.30538e-24         HES4       57801 hes family bHLH tran..
## ENSG00000187608 2.37452e-02        ISG15        9636 ISG15 ubiquitin like..
## ENSG00000188157 4.21963e-16         AGRN      375790                   agrin
## ENSG00000237330          NA       RNF223      401934 ring finger protein ..
```

Let's save our annotated data.

```r
res = res[order(res$padj),]
write.csv(res, file="deseq_results.csv")
```

## Pathway Analysis

```
library(pathview)
```

```
## ##############################################################################
## Pathview is an open source software package distributed under GNU General
## Public License version 3 (GPLv3). Details of GPLv3 is available at
## http://www.gnu.org/licenses/gpl-3.0.html. Particullary, users are required to
## formally cite the original Pathview paper (not just mention it) in publications
## or products. For details, do citation("pathview") within R.
##
## The pathview downloads and uses KEGG data. Non-academic uses may require a KEGG
## license agreement (details at http://www.kegg.jp/kegg/legal.html).
## ##############################################################################
```

```
library(gage)
```

```
##
```

```
library(gageData)

data(kegg.sets.hs)
data(sigmet.idx.hs)

#Focus on signaling and metabolic pathways only
kegg.sets.hs = kegg.sets.hs[sigmet.idx.hs]

#Examine the first 3 pathways
head(kegg.sets.hs, 3)
```

```
## $`hsa00232 Caffeine metabolism`
## [1] "10"   "1544" "1548" "1549" "1553" "7498" "9"
##
## $`hsa00983 Drug metabolism - other enzymes`
##  [1] "10"     "1066"   "10720"  "10941"  "151531" "1548"   "1549"   "1551"
##  [9] "1553"   "1576"   "1577"   "1806"   "1807"   "1890"   "221223" "2990"
## [17] "3251"   "3614"   "3615"   "3704"   "51733"  "54490"  "54575"  "54576"
## [25] "54577"  "54578"  "54579"  "54600"  "54657"  "54658"  "54659"  "54963"
## [33] "574537" "64816"  "7083"   "7084"   "7172"   "7363"   "7364"   "7365"
## [41] "7366"   "7367"   "7371"   "7372"   "7378"   "7498"   "79799"  "83549"
## [49] "8824"   "8833"   "9"      "978"
##
## $`hsa00230 Purine metabolism`
##  [1] "100"    "10201"  "10606"  "10621"  "10622"  "10623"  "107"    "10714"
##  [9] "108"    "10846"  "109"    "111"    "11128"  "11164"  "112"    "113"
## [17] "114"    "115"    "122481" "122622" "124583" "132"    "158"    "159"
## [25] "1633"   "171568" "1716"   "196883" "203"    "204"    "205"    "221823"
## [33] "2272"   "22978"  "23649"  "246721" "25885"  "2618"   "26289"  "270"
## [41] "271"    "27115"  "272"    "2766"   "2977"   "2982"   "2983"   "2984"
## [49] "2986"   "2987"   "29922"  "3000"   "30833"  "30834"  "318"    "3251"
## [57] "353"    "3614"   "3615"   "3704"   "377841" "471"    "4830"   "4831"
```

```
## [65] "4832"   "4833"   "4860"   "4881"   "4882"   "4907"   "50484"  "50940"
## [73] "51082"  "51251"  "51292"  "5136"   "5137"   "5138"   "5139"   "5140"
## [81] "5141"   "5142"   "5143"   "5144"   "5145"   "5146"   "5147"   "5148"
## [89] "5149"   "5150"   "5151"   "5152"   "5153"   "5158"   "5167"   "5169"
## [97] "51728"  "5198"   "5236"   "5313"   "5315"   "53343"  "54107"  "5422"
## [105] "5424"  "5425"   "5426"   "5427"   "5430"   "5431"   "5432"   "5433"
## [113] "5434"  "5435"   "5436"   "5437"   "5438"   "5439"   "5440"   "5441"
## [121] "5471"  "548644" "55276"  "5557"   "5558"   "55703"  "55811"  "55821"
## [129] "5631"  "5634"   "56655"  "56953"  "56985"  "57804"  "58497"  "6240"
## [137] "6241"  "64425"  "646625" "654364" "661"    "7498"   "8382"   "84172"
## [145] "84265" "84284"  "84618"  "8622"   "8654"   "87178"  "8833"   "9060"
## [153] "9061"  "93034"  "953"    "9533"   "954"    "955"    "956"    "957"
## [161] "9583"  "9615"
```

```r
#Create a named vector for the gage function
foldchanges = res$log2FoldChange
names(foldchanges) = res$entrez
head(foldchanges)
```

```
##       1266      54855      1465      51232      2034      2317
## -2.422719   3.201955 -2.313738 -2.059631 -1.888019 -1.649792
```

```r
# Get the results and examine them
keggres = gage(foldchanges, gsets=kegg.sets.hs)
attributes(keggres)
```

```
## $names
## [1] "greater" "less"    "stats"
```

```r
# Look at the first few down (less) pathways
head(keggres$less)
```

```
##                                     p.geomean stat.mean       p.val
## hsa04110 Cell cycle              8.995727e-06 -4.378644 8.995727e-06
## hsa03030 DNA replication         9.424076e-05 -3.951803 9.424076e-05
## hsa03013 RNA transport           1.375901e-03 -3.028500 1.375901e-03
## hsa03440 Homologous recombination 3.066756e-03 -2.852899 3.066756e-03
## hsa04114 Oocyte meiosis          3.784520e-03 -2.698128 3.784520e-03
## hsa00010 Glycolysis / Gluconeogenesis 8.961413e-03 -2.405398 8.961413e-03
##                                       q.val set.size         exp1
## hsa04110 Cell cycle              0.001448312      121 8.995727e-06
## hsa03030 DNA replication         0.007586381       36 9.424076e-05
## hsa03013 RNA transport           0.073840037      144 1.375901e-03
## hsa03440 Homologous recombination 0.121861535       28 3.066756e-03
## hsa04114 Oocyte meiosis          0.121861535      102 3.784520e-03
## hsa00010 Glycolysis / Gluconeogenesis 0.212222694       53 8.961413e-03
```

Let's create pathway figures.

```r
pathview(gene.data=foldchanges, pathway.id="hsa04110")
```

```
## 'select()' returned 1:1 mapping between keys and columns
```

```
## Info: Working in directory /Users/pierceford/Desktop/BGGN213/github/bggn213/Class 16 Mini Project

## Info: Writing image file hsa04110.pathview.png

# A different PDF based output of the same data
pathview(gene.data=foldchanges, pathway.id="hsa04110", kegg.native=FALSE)

## 'select()' returned 1:1 mapping between keys and columns

## Info: Working in directory /Users/pierceford/Desktop/BGGN213/github/bggn213/Class 16 Mini Project

## Info: Writing image file hsa04110.pathview.pdf

## Focus on top 5 upregulated pathways here for demo purposes only
keggrespathways <- rownames(keggres$greater)[1:5]

# Extract the 8 character long IDs part of each string
keggresids = substr(keggrespathways, start=1, stop=8)
keggresids

## [1] "hsa04640" "hsa04630" "hsa00140" "hsa04142" "hsa04330"

pathview(gene.data=foldchanges, pathway.id=keggresids, species="hsa")

## 'select()' returned 1:1 mapping between keys and columns

## Info: Working in directory /Users/pierceford/Desktop/BGGN213/github/bggn213/Class 16 Mini Project

## Info: Writing image file hsa04640.pathview.png

## 'select()' returned 1:1 mapping between keys and columns

## Info: Working in directory /Users/pierceford/Desktop/BGGN213/github/bggn213/Class 16 Mini Project

## Info: Writing image file hsa04630.pathview.png

## 'select()' returned 1:1 mapping between keys and columns

## Info: Working in directory /Users/pierceford/Desktop/BGGN213/github/bggn213/Class 16 Mini Project

## Info: Writing image file hsa00140.pathview.png

## 'select()' returned 1:1 mapping between keys and columns

## Info: Working in directory /Users/pierceford/Desktop/BGGN213/github/bggn213/Class 16 Mini Project

## Info: Writing image file hsa04142.pathview.png
```

## Info: some node width is different from others, and hence adjusted!

## 'select()' returned 1:1 mapping between keys and columns

## Info: Working in directory /Users/pierceford/Desktop/BGGN213/github/bggn213/Class 16 Mini Project

## Info: Writing image file hsa04330.pathview.png

Display images.

**Edge types**

| | |
|---|---|
| compound | ———▶ |
| hidden compound | ———▶ |
| activation | ———▶ |
| inhibition | ———⊣ |
| expression | - - - -▶ |
| repression | - - - -⊣ |
| indirect effect | ·········▶ |
| state change | ·········· |
| binding/association | – – – – |
| dissociation | ·········▶ |
| phosphorylation | —+p—▶ |
| dephosphorylation | —−p—▶ |
| glycosylation | —+g—▶ |
| ubiquitination | —+u—▶ |
| methylation | —+m—▶ |
| others/unknown | - -?- -▶ |

Thymus

Lymphoid Related
Dendritic cell

-1    0    1

IL-7

γδ T cell

CD8 T cell

CD4 T cell

Regulatory T cell

NKT cell

SCF
IL-7

SCF
IL-7

(IL-7)

Pro T cell
(DN2)

DN3

DN4

Intermediate
single-positive
cell (ISP)

Double-positive
cell (DP)

(CD2)
CD38
(CD71)
CD127
HLA-DR

(CD5)
CD25
CD44
CD117
TdT

CD5
CD7
CD38
CD71
(CD127)

CD25
CD44
CD117
TdT

CD1
(CD4)
CD7

CD2
CD5
CD38
(CD117)

CD2
CD4or8
CD5
CD7

CD3
CD5
CD38

CD2
CD4or8
CD7

CD3
CD5

| SCF | IL-7 |
|-----|------|

| HLA-DR | CD44 | CD117 | CD25 | CD127 | TdT | CD71 | CD38 | CD7 | CD2 | CD5 | CD1 | CD4 | CD8 | CD3 |
|--------|------|-------|------|-------|-----|------|------|-----|-----|-----|-----|-----|-----|-----|

NK cell Precursor

NK cell

IL-7

Lymphoid
stem cell,
Double-negative
cell (DN1)

Pro B Cell

Pre B I cell

Pre B II cell

Immature B cell

B Cell

CD34
CD44
CD117
TdT
HLA-DR

(CD9)
CD19
CD22
CD117
CD127
TdT

(CD10)
CD20
CD24
CD38
CD117
HLA-DR

CD9
CD19
CD22
CD38
CD117
TdT

CD10
CD20
CD24
CD38
HLA-DR

(CD9)
CD20
CD22
CD24
CD37
IgM

CD19
CD21
HLA-DR

(CD5)
CD19
CD21
(CD23)
CD35
HLA-DR
IgM

(CD9)
CD20
CD24
CD37
IgD

| IL-7 |
|------|

| TdT | CD117 | CD10 | CD38 | CD127 | CD9 | HLA-DR | CD19 | CD22 | CD24 | CD25 | CD20 | CD21 | CD37 | IgM | CD23 | CD35 | IgD |
|-----|-------|------|------|-------|-----|--------|------|------|------|------|------|------|------|-----|------|------|-----|

Hematopoietic
stem cell

CD34
CD135

| SCF | IL-7 |
|-----|------|

| CD34 | CD135 | TdT | HLA-DR |
|------|-------|-----|--------|

SCF
IL-3
IL-4

SCF
IL-4

Mast cell

CFU-Mast

| SCF | IL-3 | IL-4 |
|-----|------|------|

SCF
GM-CSF   IL-3

GM-CSF
IL-3

GM-CSF
IL-3

GM-CSF
IL-3

CFU-Bas

Myeloblast

Basophilic
Myelocyte

Basophil

| SCF | IL-3 | GM-CSF |
|-----|------|--------|

Flt3L
SCF   GM-CSF
IL-3

GM-CSF
IL-3
IL-5

GM-CSF
IL-3
IL-5

GM-CSF
IL-5

CFU-E0

Myeloblast

Eosinophilic
Myelocyte

Eosinophil

| Flt3L | SCF | IL-3 | GM-CSF | IL-5 |
|-------|-----|------|--------|------|

Flt3L   GM-CSF
SCF    IL-4   TNF

GM-CSF
IL-4

Myeloid Related
Dendritic Cell

Flt3L
CSF
GM-CSF

IL-3
TNF

GM-CSF
M-CSF
IL-3

GM-CSF
M-CSF
IL-3

GM-CSF
M-CSF
IL-3

GM-CSF
IL-4

CFU-M/DC

Monoblast

Promonocyte

Monocyte

GM-CSF
M-CSF

Macrophage

CD11b
CD14
CD33
CD115
CD123
CD126

CD13
CD15
CD64
CD116
CD121
CD124

CD11b
CD14
CD33
CD64
CD116
CD123
HLA-DR

CD13
CD15
CD115
CD121
CD124
CD126

CD11b
CD14
CD33
CD64

CD13
CD15
CD115
CD123
CD124
CD126

CD11b
CD14
CD33
CD64

| Flt3L | SCF | IL-3 | GM-CSF | TNF | IL-4 | M-SCF |
|-------|-----|------|--------|-----|------|-------|

| HLA-DR | CD116 | CD123 | CD33 | CD124 | CD126 | CD64 | CD115 | CD13 | CD11b | CD14 |
|--------|-------|-------|------|-------|-------|------|-------|------|-------|------|

Flt3L
SCF
G-CSF
IL-1
IL-3
IL-6
IL-11

Flt3L
SCF
GM-CSF
IL-3

GM-CSF
G-CSF
IL-3

Flt3L
SCF
GM-CSF
IL-3

GM-CSF
G-CSF

GM-CSF
G-CSF

GM-CSF
G-CSF

Myeloid
Stem Cell

CFU-GEMM

CFU-GM

CFU-G

Myeloblast

Neutrophilic
Myelocyte

Neutrophil

CD33
CD116
CD121
IL-9R
HLA-DR

CD34
CD114
CD123
EPOR

CD15
CD34
CD114
CD116
CD123
CD125
CD126

CD33
CD64
CD115
CD121
CD124
HLA-DR

CD13
CD33
CD116
CD123
CD125
HLA-DR

CD15
CD114
CD115
CD121
CD124
CD126

CD13
CD33
CD116
CD123
CD125

CD15
CD114
CD124
CD126

CD11b
CD15
CD116
CD125

CD11b
CD15
CD33

Bone marrow

| Flt3L | SCF | G-SCF | IL-3 | IL-6 | IL-11 | IL-1 | GM-CSF |
|-------|-----|-------|------|------|-------|------|--------|

| Flt3L | SCF | IL-3 | GM-CSF | G-CSF |
|-------|-----|------|--------|-------|

| IL-9R | CD34 | HLA-DR | CD116 | CD121 | CD114 | CD123 | CD124 | CD126 | CD33 | CD13 | CD125 | CD11b |
|-------|------|--------|-------|-------|-------|-------|-------|-------|------|------|-------|-------|

Flt3L
SCF
GM-CSF
IL-4

IL-3
IL-4

SCF
GM-CSF
IL-4   EPO

TPO
EPO

EPO

BFU-E

CFU-E

Proerythroblast

Erythrocyte

CD33
CD117
EPOR

CD34
CD123
HLA-DR

CD36
CD235a

CD235a

CD35
CD55
CD235a

CD44
CD59

| Flt3L | SCF | GM-CSF | IL-3 | IL-4 | EPO | TPO |
|-------|-----|--------|------|------|-----|-----|

| HLA-DR | EPOR | CD33 | CD34 | CD117 | CD123 | CD36 | CD235a | CD35 | CD44 | CD55 | CD59 |
|--------|------|------|------|-------|-------|------|--------|------|------|------|------|

Flt3L
SCF
GM-CSF
IL-3

IL-6
IL-11
TPO

Flt3L
SCF
GM-CSF

Meg-CSF
IL-3
IL-6

IL-11
TPO

SCF
GM-CSF
IL-3

IL-6
IL-11
TPO

IL-6
IL-11
TPO

BFU-MK

CFU-MK

Mega-
karyocyte

Platelets

CD33
CD116
CD126
HLA-DR

CD34
CD123
IL-11R

CD61
CD116
CD122
CD126

CD9
CD36
CD42
CD116
CD126

CD14
CD41
CD61
CD123

CD9
CD36
CD42
CD61

CD14
CD49
CD126

| Flt3L | SCF | IL-3 | IL-6 | IL-11 | GM-CSF | Meg-CSF | TPO |
|-------|-----|------|------|-------|--------|---------|-----|

| HLA-DR | CD33 | CD34 | IL-11R | CD116 | CD123 | CD126 | CD61 | CD9 | CD14 | CD36 | CD41 | CD42 | CD49 |
|--------|------|------|--------|-------|-------|-------|------|-----|------|------|------|------|------|

Data on KEGG graph
Rendered by Pathview

JAK-STAT SIGNALING PATHWAY

-1    0    1

Cytokine-cytokine
receptor interaction

ECS complex

Ubiquitin
mediated proteolysis

STAM

Cytokine

Hormone

GF

Receptor  JAK  +p

-p    -p

TC-PTP  SHP1

STAT

STAT
STAT

STAT dimerization

IRF9

TC-PTP  PIAS

-p

CBP/P300  SLIM

+u

Proteasome

DNA

CIS  SOCS

Bcl-2  MCL1

Bcl-XL  PIM1

c-Myc  CycD

p21

AOX

GFAP

Anti-apoptosis

Cell-cycle progression

Cell-cycle inhibition

Lipid metabolism

Differentiation

Apoptosis

Cell cycle

+p  SHP2
GRB

SOS

Ras

MAPK
signaling pathway

Raf

Proliferation
Differentiation

+p  PI3K

AKT

mTOR

PI3K-AKT
signaling pathway

Cell cycle
Cell survival

Data on KEGG graph
Rendered by Pathview

16

STEROID HORMONE BIOSYNTHESIS

Steroid biosynthesis

Cholesterol sulfate

Cholesterol

20α-Hydroxy-cholesterol
22β-Hydroxy-cholesterol
20α,22β-Dihydroxy-cholesterol
17α,20α-Dihydroxy-cholesterol

4-Methylpentanal

Pregnenolone
21-Hydroxy-pregnenolone
7α-Hydroxy-pregnenolone
17α-Hydroxy-pregnenolone
17α,21-Dihydroxy-pregnenolone
11β,17α,21-Trihydroxy-pregnenolone

Progesterone
11α-Hydroxy-progesterone
20α-Hydroxy-progesterone
11β-Hydroxy-progesterone
11-Deoxy-corticosterone
17α-Hydroxy-progesterone
17α,20α-Dihydroxy-pregn-4-en-3-one
11-Deoxycortisol

Corticosterone
18-Hydroxy-corticosterone
Aldosterone hemiacetal
CYP11B2
Aldosterone
11-Dehydro-corticosterone
21-Hydroxy-5β-pregnane-3,11,20-trione
Tetrahydro-corticosterone

5α-Dihydro-deoxycorticosterone
Allotetrahydro-deoxycorticosterone

11β,21-Dihydroxy-3,20-oxo-5β-pregnan-18-al
3α,11β,21-Trihydroxy-20-oxo-5β-pregnan-18-al
11β,21-Dihydroxy-5β-pregnane-3,20-dione
3α,20α,21-Trihydroxy-5β-pregnane-11-one
3α,21-Dihydroxy-5β-pregnane-11,20-dione

C21-Steroids

5β-Pregnane-3,20-dione
3α-Hydroxy-5β-pregnan-20-one
Pregnanediol
5α-Pregnane-3,20-dione
3α-Hydroxy-5α-pregnan-20-one
5α-Pregnan-20α-ol-3-one
5α-Pregnane-3α,20α-diol

21-Deoxycortisol
Cortisol
4-Androsten-11beta-ol-3,17-dione
11β-Hydroxytestosterone
Urocortisol
Cortol
11β,17α,21-Trihydroxy-5β-pregnane-3,20-dione
Cortisone
17α,21-Dihydroxy-5β-pregnane-3,11,20-trione
Cortolone
Urocortisone

Dehydroepiandrosterone
Dehydroepiandrosteron sulfate
7α-Hydroxydehydro-epiandrosterone
Adrenosterone
16-Hydroxyandrost-4-ene-3,17-dione
16α-Hydroxydehydro-epiandrosterone
3β,17β-Dihydroxy-androst-5-ene

Androst-4-ene-3,17-dione
7α-Hydroxy-androstenedione
7α-Hydroxy-testosterone
Testosterone

5β-Androstane-3,17-dione
Etiocholan-3α-ol-17-one
Etiocholan-3α-ol-17-one 3-glucuronide
5α-Androstane-3,17-dione
Androsterone
Androsterone-glucuronide
19-Hydroxyandrost-4-ene-3,17-dione
19-Oxoandrost-4-ene-3,17-dione
3-Oxo-13,17-secoandrost-4-ene-17,13a-lactone
19-Hydroxy-testosterone
19-Oxotestosterone
5α-Dihydro-testosterone
Androstan-3alpha,17beta-diol
5β-Dihydro-testosterone
Testosterone glucuronide

C19-Steroids

Estrone 3-sulfate
Estrone glucuronide
Estradiol-17α
2-Methoxyestrone-3-glucuronide
2-Methoxyestrone
2-Hydroxyestrone
2-Methoxyestrone-3-sulfate
16-α-Hydroxyestrone
Estrone
Estradiol-17β
16-Glucuronide-estriol
Estriol
2-Methoxy-estradiol-17β-3-glucuronide
2-Methoxy-estradiol-17β
2-Hydroxy-estradiol-17β
2-Methoxy-estradiol-17β-3-sulfate
6β-Hydroxy-estradiol-17β
Estradiol-17β-3-glucuronide
Estradiol-17β-3-sulfate

C18-Steroids

Data on KEGG graph
Rendered by Pathview

-1   0   1

LYSOSOME

-1  0  1

bacterium

cytosol
pH~ 7.2

lysosomal acid hydrolase

Phagocytosis

Golgi body

ATP    ADP

transport vesicle

phagosome

ATPa-V

H+

pH~ 5.0

clathrin coat

Transport of
synthesized lysosomal enzymes
(See below)

Endocytosis

Endocytosis

early endosome

late endosome
multivesicular body (MVB)

acid hydrolase

lysosomal membrane protein

MCOLN1

lysosome

mitochondria

autophagosome

Glycosaminoglycan
degradation

Other glycan
degradation

plasma membrane

Autophagy

Regulation of autophagy

Lysosomal acid hydrolases
proteases
| cathepsin | napsin | LGMN | TPP1 |

glycosidases
| GLA | GLB | GAA | GBA | IDUA |
| NAGA | NAGLU | GALC | GUSB | FUCA1 |
| HEXA/B | MANB | LAMAN | NEU1 | HYAL1 |

sulfatases
| ARS | GALNS | GNS | IDS | SGSH |

lipases                nuclease  phosphatase
| LIPA | LYPLA3 |   | DNaseII | ACP2 | ACP5 |

sphingomyelinase   ceramidase   aspartylglucosaminidase
| SMPD1 |          | ASAH1 |      | AGA |

Other lysosomal enzymes and activators
| saposin | GM2A | CLN1 |

Lysosomal membrane proteins
major lysosomal membrane proteins
| LAMP | LIMP |

minor lysosomal membrane proteins
| NPC | cystscan | sialin | NRAMP | LAPTM |
| ABCA2 | ABCB9 | ACP2 | endolyn | LALP70 |
| sortilin | CLN3 | CLN5 | CLN7 | HGSNAT |
| MCOLN1 | LITAF |

Activation of
lysosomal sulfatase
precursor

lysosomal hydrase
precursor

M6P receptor

Receptor-dependent
transport

FGE

from ER

mannose

M6P

MPR

M6P

ATPa-V

GNPT
NAGPA

+pO

| AP-1 | AP-3 |
| GGAs | AP-4 |

clathrin

M6P

transport vesicle

AP-3

lysosome

Snare interactions
in vesicular transport

M6P

mannose

cis Golgi
network

trans Golgi
network

AP-1

Golgi body

Receptor recycling

late endosome

mature
lysosomal hydrase

Transport of synthesized lysosomal enzymes

**Data on KEGG graph**
**Rendered by Pathview**

-1  0  1

NOTCH SIGNALING PATHWAY

Fringe

Dvl

Numb

Co-activator

MAML
HATs

SKIP

Delta

Notch

Deltex

NICD

CSL

DNA

Hes1/5

Serrate

(Notch intracellular
domain)

Hey

S3

S2

PSE2  PSEN
NCSTN  APH-1

γ-Secretase complex

Hairless

SMRT

PreTα

CtBP  Gro/TLE  CIR

Co-repressor

HDAC  ATXN/1L

TACE

Ras/MAPK

MAPK signaling
pathway

Gene expression

**Data on KEGG graph**
**Rendered by Pathview**

18